

A Transport-Friendly NIC for Multicore/Multiprocessor Systems

Wenji Wu, Phil DeMar, Matt Crawford
Fermilab, P.O. Box 500, Batavia, IL 60510

Abstract - Receive side scaling (RSS) is a NIC technology that provides the benefits of parallel receive processing in multiprocessing environments. However, RSS lacks a critical data steering mechanism that would automatically steer incoming network data to the same core on which its application thread resides. This absence causes inefficient cache usage if an application thread is not running on the core on which RSS has scheduled the received traffic to be processed and results in degraded performance. To remedy the RSS limitation, Intel's Ethernet Flow Director technology has been introduced. However, our analysis shows that Flow Director can cause significant packet reordering. Packet reordering causes various negative impacts in high-speed networks. We propose a NIC data steering mechanism to remedy the RSS and Flow Director limitations. This data steering mechanism is mainly targeted at TCP. We term a NIC with such a data steering mechanism "A Transport Friendly NIC" (A-TFN). Experimental results have proven the effectiveness of A-TFN in accelerating TCP/IP performance.

Indexed Terms - TCP/IP, Parallel Network Stacks, Core Affinity, High Performance Networking, 40GigE, 100GigE.

1. Introduction & Motivation

Computing is shifting towards multiprocessing. The fundamental goal of multiprocessing is improved performance through the introduction of additional hardware threads, CPUs, or cores (all of which will be referred to as "cores" for simplicity). The emergence of multiprocessing has brought both opportunities and challenges for TCP/IP performance optimization in such environments. Modern network stacks can exploit parallel cores to allow either message-based parallelism or connection-based parallelism as a means of enhancing performance [1]. To date, major network stacks such as Windows, Solaris, and Linux have been redesigned and parallelized to better utilize additional cores. While existing OSes exploit parallelism by allowing multiple threads to carry out network operations concurrently in the kernel, supporting this parallelism carries significant costs, particularly in the context of contention for shared resources, software synchronization, and poor cache efficiencies. However, investigations [2][3][4] indicate that CPU core affinity on network processing in multiprocessing environment can significantly reduce contention for shared resources, minimize software synchronization overheads, and enhance cache efficiency.

Core affinity on networking processing has the following goals: (1) *Interrupt affinity*: Network interrupts of the same type should be directed to a single core. Redistributing network interrupts in either a random or round-robin fashion to different cores has undesirable side effects [3]. (2) *Flow affinity*: Packets belong to a specific TCP flow should be processed by the same core. TCP has a large and frequently accessed state that must be shared and protected when packets from the same connection are processed. Flow affinity reduces contention for shared resources, minimizes software

synchronization, and enhances cache efficiency. (3) *Network data affinity*: Incoming network data should be steered to the same core on which its application thread resides. This is becoming more important with the advent of Direct Cache Access (DCA) [5].

The emergence of parallel network stacks and the necessity of core affinity on network processing in multiprocessing environment require new NIC designs. An NIC should not only provide mechanisms to allow parallel receive processing to better utilize parallel network stacks, but also to facilitate core affinity on network processing in multiprocessing environments. RSS [6] is a NIC technology that steps toward that direction. RSS supports multiple receive queues; it assigns packets of the same data flow to a single queue and evenly distributes traffic flows across queues. With Message Signal Interrupt (MSI/MSI-X) support, each receive queue is assigned a dedicated interrupt and RSS steers interrupts on a per-queue basis. RSS provides the benefits of parallel receive processing in multiprocessing environments. However, RSS has a limitation: it cannot steer incoming network data to the same core where its network application thread resides. The reason is simple: RSS does not maintain the relationship "Traffic Flows \rightarrow Network applications \rightarrow Cores" in the NIC. Since network applications run on cores, we simply put it as "Traffic Flows \rightarrow Cores (Applications)". This is symptomatic of a broader disconnect between existing software architecture and multicore hardware. With OSes like Windows and Linux, if an application thread is running on one core, while RSS has scheduled received traffic to be processed on a different core, poor cache efficiency and significant core-to-core synchronization overheads will result. The overall system efficiency may be severely degraded (see Section 2).

In parallel to our research, Intel has introduced the Ethernet Flow Director technology [7]. The basic idea is simple: Flow Director maintains the relationship "Traffic Flows \rightarrow Cores (Applications)" in the NIC. Flow Director not only provides the benefits of parallel receive processing in multiprocessing environments, it also can automatically steer packets of a specific data flow to the same core on which its application thread resides. However, our research shows that Flow Director can cause significant packet reordering in multiprocessing environments [see Section 2.5]. In high-speed networks, packet reordering causes various negative impacts [8][9]. In addition, TCP Selective Acknowledgement (SACK) is now implemented and enabled by almost all general-purpose OSes. When packet reordering occurs, processing or generating TCP SACK information can seriously degrade the TCP sender or receiver's performance [10]. For example, the receiver would sort the out-of-order queue to generate SACKs in the event of packet reordering. Sorting the out-of-order queue is expensive, especially when the queue is large. Because the networking community is working towards

40GigE and 100GigE, the performance requirements on TCP/IP are becoming more challenging. Flow Director's packet reordering problem becomes more serious.

We propose a NIC mechanism to remedy the RSS and Flow Director limitations. It steers incoming network data to the same core on which its application thread resides and ensures in-order packet delivery. Our data steering mechanism is mainly targeted at TCP, but can be extended to UDP and SCTP. We term a NIC with such a data steering mechanism A Transport-Friendly NIC, or A-TFN. As Flow Director, A-TFN maintains the relationship "Traffic Flows \rightarrow Cores (Applications)" in the NIC, with OSes correspondingly enhanced to support such capability. For transport layer traffic, A-TFN maintains a Flow-to-Core table in the NIC, with one entry per flow. Each entry tracks which core a flow should be assigned to. However, A-TFN is different from Flow Director in two significant ways: (1) A-TFN applies a very simple yet effective mechanism to update the Flow-to-Core table in the NIC. It requires the most minimal OS support (see Section 3). However, to support Flow Director, OS must be multiple TX queue capable [11]. Therefore, A-TFN is simpler and minimizes changes in the OS. And (2) A-TFN has a mechanism to ensure in-order packet delivery. Flow Director does not have such a mechanism and our analysis and experiments show that Flow Director can cause significant packet reordering in multiprocessing environments.

To design A-TFN, there is an obvious trade-off between the amount of work done in the NIC and in the OS. In the paper, we discuss two design options. Option 1 is to minimize changes in the OS and focuses instead on identifying the minimal set of mechanisms to add to the NIC. This design adds complexity and cost to the NIC. On the other end of the design space, it could be let the OS update the flow-to-core table directly without changing anything in the NIC hardware (option 2). Conceptually, this approach could be fairly straightforward to implement. However, it might add significant extra communication overheads between the OS and the NIC, especially when the Flow-to-Core table gets large. Due to space limitation, this paper is mainly focused on the first design option. The new NIC is emulated in software and it shows that the solution is effective and practical to remedy the limitations in RSS and Flow Director. In future work, we will explore the second design option.

The contributions of this paper are fourfold. First, we show for certain OSes, such as Linux, that tying a traffic flow to a single core does not necessarily ensure flow affinity or network data affinity. Second, we show that RSS lacks a mechanism to automatically steer packets of a data flow to the same core(s) on which its application thread resides. Third, we show that Flow Director can cause significant packet reordering in multiprocessing environments. Flow Director lacks mechanisms to ensure in-order packet delivery when it steers packets across cores. Fourth, we propose the A-TFN mechanism to remedy the limitations in RSS and Flow Director. Experimental results have proven the effectiveness of A-TFN in accelerating TCP/IP performance. Because the networking community is headed for 40GigE and 100GigE,

the performance requirements on TCP/IP are more challenging; further architecture optimizations and technology advances are necessary. A-TFN is working toward that direction and is timely.

The remainder of the paper is organized as follows: In Section 2, we present problem formulation. Section 3 describes the A-TFN mechanism. In section 4, we discuss experiment results that showcase the effectiveness of our A-TFN mechanism. In section 5, we present related research. We conclude in section 6.

2. Problem Formulation

2.1 Packet Receive Processing with RSS

RSS is a NIC technology. It supports multiple receive queues and integrates a hashing function in the NIC. NIC computes a hash value for each incoming packet. Based on hash values and an indirection table, NIC assigns packets of the same data flow to a single queue and evenly distributes traffic flows across queues. With Message Signal Interrupt (MSI/MSI-X) and Flow Pinning support, each receive queue is assigned a dedicated interrupt and tied to a specific core. The device driver allocates and maintains a ring buffer for each receive queue within system memory. For packet reception, a ring buffer must be initialized and pre-allocated with empty packet buffers. The ring buffer size is device- and driver-dependent. Fig. 1 illustrates packet receive-processing with RSS: (1) When incoming packets arrive, the hash function is applied to the header to produce a hash result. The hash result is used to index the indirection table. The indirection table is the data structure that contains an array of core numbers to be used for RSS. Each lookup from the indirection table identifies the core and hence, the associated receive queue. (2) The NIC assigns incoming packets to the corresponding receive queues. (3) The NIC DMAs (direct memory access) the received packets into the corresponding ring buffers in the host system memory. (4) The NIC sends interrupts to the cores that are associated with the non-empty queues. Subsequently, the cores respond to the network interrupts and process received packets up through the network stack from the corresponding ring buffers one by one. Appendix A has more details of RSS mechanisms.

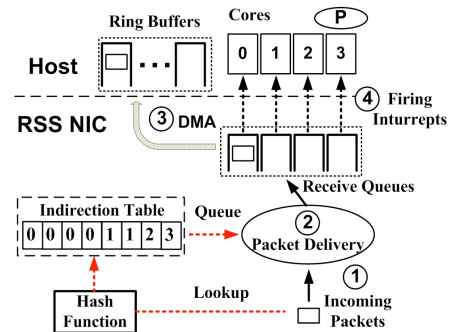


Fig. 1 Packet Receiving Process with RSS

The OS can periodically rebalance the network load on cores by updating the indirection table, based on the assumption that the hash function will evenly distribute

incoming traffic flows across the indirection table entries. Since the OS does not know which specific entry in the indirection table an incoming traffic flow will be mapped to, it can only passively react to load imbalance situations by changing each core's number of appearances in the indirection table. For better load balancing performance, the size of the indirection table is typically two to eight times the number of cores in the system [6]. For example, in Fig. 1, the indirection table has 8 entries, which are populated as shown. As such, traffic loads directed to Core 0, 1, 2, and 3 are 50%, 25%, 12.5%, and 12.5%, respectively.

2.2 RSS Limitation and the Reasons

RSS provides the benefits of parallel receive processing. However, this mechanism does present certain limitation: it cannot steer incoming network data to the same core on which its application thread resides. The reason is simple: RSS does not maintain the relationship “Traffic Flows → Cores (Applications)” in the NIC. When packets arrive, the hash function is applied to the header to produce a hash result. Based on the hash values, the NIC assigns packets to receive queues and then cores, with no way to consider on which core the corresponding application thread is running. Although receive queues can be instructed to send interrupt to a specific set of cores, existing general purpose OSes can only provide limited process-to-interrupt affinity capability; network interrupt delivery is not synchronized with process scheduling. This is because the OS schedulers have other priorities, such as load balancing and fairness, over process-to-interrupt affinity. Besides, multiple network applications' traffic might map to a single interrupt, which brings new challenges to an OS scheduler. Therefore, a network application thread might be scheduled on cores other than those where its corresponding network interrupts are directed. This is symptomatic of a broader disconnect between existing software architecture and multicore hardware.

OSes like Windows implement the function of the indirection table, which can provide limited data steering capabilities for RSS. However, it still cannot steer packets of a data flow to the same core where the application thread resides. Turning again to Fig 1, network application thread P is scheduled to run on Core 3. Its traffic might be hashed to an entry that directs to other cores. The OS does not know which specific entry in the indirection table a traffic flow will be mapped to.

With existing RSS capability, there are many cases in OSes in which a network application resides on cores other than those to which its corresponding network interrupts are directed: (1) A single-threaded application might handle multiple concurrent TCP connections. Assuming such an application handles n concurrent TCP connections and runs on an m -core system, an RSS-enabled NIC will evenly (statistically) distribute the n connections across the m cores. Since the application thread can only run on a single core at any moment, only n/m connections' network interrupts are directed to the same core where the application runs. (2) Soft partition technologies like CPUSSET [12] are applied in the

context of networking environments. Since the OS (or system administrator) has no way of knowing to which specific core they will be mapped, network applications might be soft-partitioned on cores other than those to which their network interrupts are directed. (3) The general purpose OSes scheduler prioritizes load balancing or power saving over process-to-interrupt affinity [13][14]. For OSes like Linux, when the multicore peak performance mode is enabled, the scheduler tries to use all cores in parallel to the greatest extent possible, distributing the load equally among them. When the multicore power saving mode is enabled, the scheduler is biased to restrict the workload to a single physical processor. As a result, a network application might be scheduled on cores other than those to which its network interrupts are directed. For clarity, we illustrate the above cases in Fig. 2. The system contains two physical processors, each with two cores. P1 – P5 are processes that run within the system. P1 is a network application thread that includes traffic flows. An RSS-enabled NIC steers the traffic flows to different cores, as shown in the figure (red arrows). In all of these cases, P1 resides on cores other than those to which its corresponding network interrupts are directed.

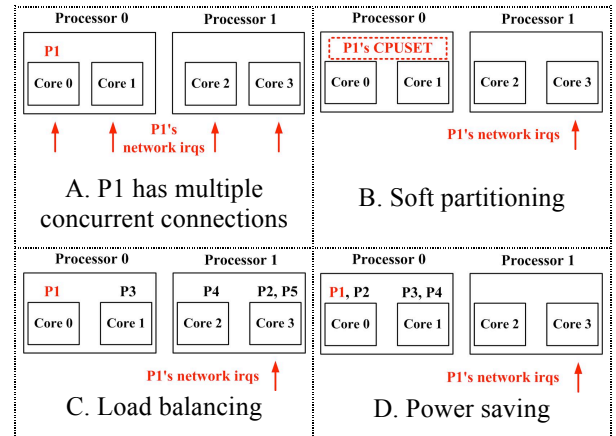


Fig. 2 Network Irqs and Apps. on Different Cores

On OSes like Windows, when a core responds to the network interrupt, the corresponding interrupt handler is called, within which a deferred procedure call (DPC) is scheduled. On the core, DPC processes received packets up through the network stack from the corresponding ring buffer one by one [15]. Therefore, on Windows, tying a traffic flow to a single core does ensure interrupt affinity and flow affinity. However, if network interrupts are not directed to cores on which the corresponding applications reside, network data affinity cannot be achieved, resulting in degraded cache efficiency [6]. This reality might cause serious performance degradation for NUMA systems. On some OSes, like Linux, tying a traffic flow to a single core does not necessarily ensure flow affinity or network data affinity due to Linux TCP's unique prequeue-backlog queue design. In the following sections, we discuss in detail why the combination of RSS and Flow Pinning cannot ensure flow affinity and network data affinity in Linux.

2.3 Linux Network Processing in Multicore Systems

Linux allows multiple threads to simultaneously process different packets from the same or different connections. Two types of threads may perform network processing in Linux: application threads in process context and interrupt threads in interrupt context. When an application makes socket-related system calls, that application's process context may be borrowed to carry out network processing. When a NIC interrupts a core, the associated handler services the NIC and schedules the softirq, softnet. Afterwards, the softnet handler processes received packets up through the network stack in interrupt context. TCP is a connection-oriented protocol, and it has a large and frequently accessed state that must be shared and protected. In the case of the Linux TCP, the data structure socket maintains a connection's various TCP states, and there is a per-socket lock to protect it from unsynchronized access. The lock consists of a spinlock and a binary semaphore. The binary semaphore construction is based on the spinlock. In Linux, since an interrupt thread cannot sleep, when it accesses a socket, the socket is protected with the spinlock. When an application thread accesses a socket, the socket is locked with the binary semaphore and is considered "owned-by-user." The binary semaphore synchronizes multiple application threads among themselves. It is also used as a flag to notify interrupt threads that a socket is "owned-by-user" to coordinate synchronized access to the socket between interrupt and application threads. Our previous research [16][17] studied the details of the Linux packet receiving process. Here, we simply summarize Linux TCP processing of the data receive path in interrupt and process contexts, respectively.

a) TCP Processing in Interrupt Context

- (1) When the NIC interrupts a core, the network interrupt's associated handler services the NIC and schedules the softirq, softnet.
- (2) The softnet handler moves a packet from the ring buffer and processes the packet up through the network stack. If there is no packet available in the ring buffer, the softnet handler exits.
- (3) A TCP packet (segment) is delivered up to the TCP layer. The network stack first tries to identify the socket to which the packet belongs, and then seeks to lock (spinlock) the socket.
- (4) The network stack checks if the socket is "owned-by-user" or if an application thread is sleeping and awaiting data:
 - If yes, the packet will be enqueued into the socket's backlog queue or prequeue. TCP processing will be performed later in process context by the application thread.
 - If not, the network stack will perform TCP processing on the packet in interrupt context.
- (5) Unlock the socket; go to step 2.

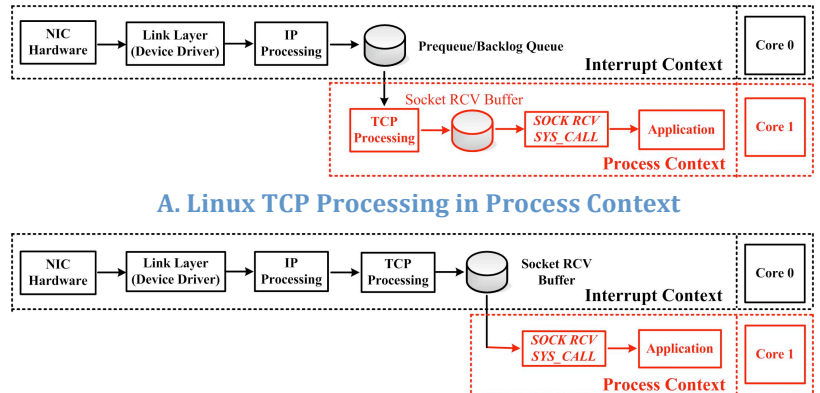
b) TCP Processing in Process Context

- (1) An application thread makes a socket-related receive system call.

- (2) Once the system call reaches the TCP layer, the network stack seeks to lock (semaphore) the socket first.
- (3) The network stack moves data from the socket into the user space, and generates ACKs.
- (4) If the socket's prequeue and/or backlog queue are not empty, the calling application's process context will be borrowed to carry out TCP processing.
- (5) Unlock the socket and return from the system call.

For the data transmit path, network processing starts in the process context when an application makes socket-related system calls to send data. If TCP gives permission to send (based on TCP receiver window, congestion window, and sender window statuses), network processing in process context can reach down to the bottom of the protocol stack. Otherwise, transmit side network processing is triggered by incoming TCP ACKs for the data receive path, which are performed in their execution environments (interrupt or process contexts). In this paper, we focus mainly on receive side processing because it is known to be more memory intensive and complex, and TCP processing on the transmit side is also dependent on ACKs in the data receive path.

As described above, whether TCP processing is performed in process or interrupt contexts depends on the volatile runtime environments. For example, we used FTP to download Linux kernels from www.kernel.org and instrumented the Linux network stack to record the percentage of traffic processed in process context. The recorded percentage ranged from 50% to 75%. In a multicore system, when an application's process context is borrowed to execute the network stack, TCP processing is performed on the core(s) where the application is scheduled to run. When TCP processing is performed in interrupt context, it is performed on the cores to which the network interrupts are directed. Take, for example, Fig. 3, in which network interrupts are directed to core 0 and the associated network application thread is scheduled to run on core 1. In interrupt context, TCP is processed on core 0; in process context, this occurs on core 1. Since TCP processing performed in process or interrupt contexts depends on volatile runtime conditions, it may alternate between these two cores. Therefore, although the combination of RSS and Flow Pinning can tie a traffic flow to a single core, when a network application thread resides on some other core, TCP processing might alternate between



B. Linux TCP Processing in Interrupt Context

Fig. 3 Linux TCP Processing Contexts in the Data Receive Path

different cores. We would achieve neither flow affinity nor network data affinity.

2.4 Negative Impacts

If a network application runs on cores other than those where its corresponding RSS network interrupts are directed, various negative impacts result. On both Windows and Linux systems, network data affinity cannot be achieved. On OSes like Linux, TCP processing might alternate between different cores even if the interrupts for the flow are pinned to a specific core. As a result, it will lead to poor cache efficiency and cause significant core-to-core synchronization overheads. Also, it renders the DCA technology ineffective. In multiple core systems, core-to-core synchronizations involve costly snoops and MESI operations [18], resulting in extra system bus traffic. This is especially expensive when the contending cores exist within different physical processors, which usually involves synchronous read/write operations to a certain memory location. In addition, for Linux, interrupt and application threads contend for shared resources, such as locks, when they concurrently process packets from the same flow. The socket's spinlock, for example, would be in severe contention. When a lock is in contention, contending threads simply wait in a loop ("spin"), repeatedly checking until the lock becomes available. While waiting, no useful work is executed. Contention for other shared resources, such as memory and system bus, also occurs frequently. Since this intra-flow contention may occur on a per-packet basis, the total contention overhead could be severe in high network I/O environments.

To demonstrate the negative impacts, we ran data transmission experiments over an isolated sub-network.

Sender: Dell R-805; 2 Quad Core AMD Opteron 2346HE, 1.8GHz; Broadcom NetXtreme II 1Gbps NIC; Linux 2.6.28. **Receiver:** SuperMicro Server; 2 Intel Xeon CPUs, 2.66 GHz; Intel PRO/1000 1Gbps NIC (DCA not supported); Linux 2.6.28. The receiver's CPU architecture is as shown in Fig. 4.

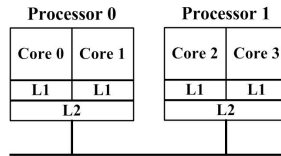


Fig. 4 Receiver CPUs

In the experiments we used iperf [19] to send data in one direction. The sender transmitted one TCP stream to the receiver for 100 seconds. In the receiver, network interrupts were all directed to core 0. However, iperf was pinned to different cores: (1) Iperf was pinned to core 0 (network interrupts and applications were pinned to the same core). (2) Iperf was pinned to core 1 (network interrupts and applications were pinned to different cores, but within the same processor). (3) Iperf was pinned to core 2 (network interrupts and applications were pinned to different processors). The throughput rates in these experiments all saturated the 1Gbps link (around 940 Mbps). The experiments were designed to feature the same throughput rates. Therefore, we do not need to normalize the final results with the throughputs. We ran *oprofile* [20] to profile system performance in the case of the receiver. The metrics of interest were: INST_RETIRED, the

number of instructions retired; BUS_TRAN_ANY, the total number of completed bus transactions; and BUS_HITM_DRV, the number of HITM (hit modified cache line) signals asserted [21]. For these metrics, the number of events between samples was 10000. We also enabled the Linux Lockstat [12] to collect lock statistics. On this basis we calculated the total time spent waiting to acquire various kernel locks, and we called this WAITTIME-TOTAL. Consistent results were obtained across repeated runs. The results are as listed in Fig. 5, with a 95% confidence interval.

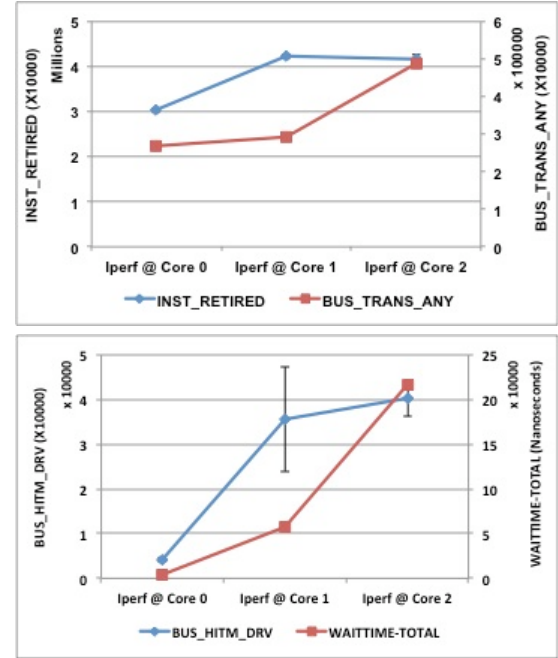


Fig. 5. Experiment Results

The throughput rates in these experiments all saturated the 1Gbps link. However, Fig. 5 shows that the metrics of iperf @ Core 1 and Core 2 are much higher than those of iperf @ Core 0. This verifies that when a network application is scheduled on cores other than those to which the corresponding system interrupts are directed, severely degraded system efficiency will result. INST_RETIRED measures the load on the receiver. The results demonstrate that contention for shared resources between interrupt and application threads led to an extra load. The extra load is mainly related to time spent waiting for locks. The experimental WAITTIME-TOTAL data verify this point. It is surprising that the BUS_TRAN_ANY of iperf @ Core 2 is almost twice that of iperf @ Core 0. The BUS_HITM_DRV of iperf @ Core 0 is far less than that of iperf @ Core 1 and Core 2. Since the throughput rates in these experiments all saturated the 1Gbps link, the extra BUS_TRAN_ANY and BUS_HITM_DRV transactions of iperf @ Core 1 and Core 2 were caused by cache trashing and lock contention, as analyzed above.

2.5 Why does Flow Director cause packet reordering?

Intel has introduced the Ethernet Flow Director technology to remedy the RSS limitation. Flow Director is a NIC technology. As shown in Fig. 6, it supports multiple

receive queues in the NIC, up to the number of cores in the system. Each receive queue has a dedicated interrupt and is tied to a specific core; each core in the system is assigned a specific receive queue. Flow Director maintains

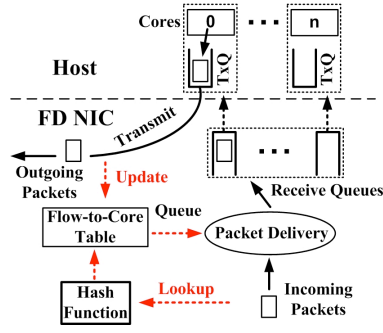


Fig. 6 Flow Director Mechanism

a Flow-to-Core table with a single entry per flow. Each entry tracks the receive queue (core) to which a flow should be assigned. Entries within the Flow-to-Core table are updated by outgoing packets. To support Flow Director, OS must be multiple TX queue capable [11]. Each core in the system is assigned a specific transmit queue. Outgoing traffic generated on a specific core is transmitted via its corresponding transmit queue. For an outgoing transport-layer packet, the OS records the processing core ID and use it to update the corresponding entry in the table. Flow Director makes use of the 5-tuple $\{src_addr, dst_addr, protocol, src_port, dst_port\}$ in the receive direction to specify a flow. Therefore, for an outgoing packet with the header $\{(src_addr: x), (dst_addr: y), (protocol: z), (src_port: p), (dst_port: q)\}$, its corresponding flow entry in the table is identified as $\{(src_addr: y), (dst_addr: x), (protocol: z), (src_port: q), (dst_port: p)\}$. Packet receiving process with Flow Director is similar to that of with RSS, except that incoming packets look up the Flow-to-Core table to identify the core.

Flow Director not only provides the benefits of parallel receive processing in multiprocessing environments, it also can automatically steer packets of a data flow to the same core on which its application resides. However, our analysis shows that Flow Director cannot guarantee in-order packet delivery in multiprocessing environments.

As shown in Fig. 7, at time $T - \epsilon$, Flow 1's flow entry maps to Core 0 in the Flow-to-Core table. At this instant, packet S of Flow 1 arrives; based on the "Traffic Flow \rightarrow Core" table, it is assigned to Core 0. At time T , due to process migration, Flow 1's flow entry is updated and maps to Core 1. At $T + \epsilon$, Packet $S+1$ of Flow 1 arrives and is assigned to the new core, namely Core 1. After assigning received packets to the corresponding receive queues, NIC copies them into system memory via DMA, and fires network interrupts, if necessary. When a core responds to a network interrupt, it processes received packets up through the network stack from the corresponding ring buffer one by one. In

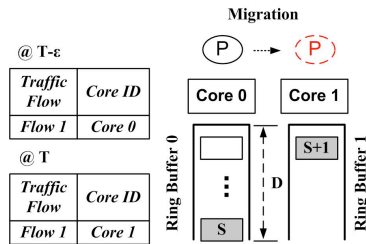


Fig. 7 A Simplified Model for Packet Reordering Analysis

our case, Core 0 processes packet S up through the network stack from Ring Buffer 0, and Core 1 services packet $S+1$ from Ring Buffer 1. Let $T_{service}(S)$ and $T_{service}(S+1)$ be the times at which the network stack starts to service packets S and $S+1$, respectively. If $T_{service}(S) > T_{service}(S+1)$, the network stack would receive packet $S+1$ earlier than packet S , resulting in packet reordering. Let D be the ring buffer size and let the network stack's packet service rate be $R_{service}$ (packets per second). Assume there are n packets ahead of S in Ring Buffer 0 and m packets ahead of $S+1$ in Ring Buffer 1. Then it has

$$T_{service}(S) = T - \epsilon + n / R_{service} \quad \text{and} \quad T_{service}(S+1) = T + \epsilon + m / R_{service}.$$

If ϵ is small and $n > m$, the condition of $T_{service}(S) > T_{service}(S+1)$ would easily hold and lead to packet reordering. Since the ring buffer size is D , the worst case is $n = D - 1$ and $m = 0$. It has $T_{service}(S) = T - \epsilon + (D - 1) / R_{service}$ and $T_{service}(S+1) = T + \epsilon$. The ring buffer size D is a design parameter for the NIC and driver. For example, the Myricom 10Gb NIC is 512.

In a multicore system, a general-purpose OS scheduler tries to use all core resources in parallel as much as possible, distributing and adjusting the load among the cores. Process migration across cores occurs frequently. The conditions for Flow Director to cause packet reordering can be easily satisfied. Flow Director can easily cause packet reordering.

To validate our analysis, we ran data transmission experiments over an isolated network. A sender was directly connected to a receiver via a physical 10Gbps link. The sender and receiver are the same computer systems as specified in Section 2.4, except that:

Sender: Myricom 10Gbps Ethernet NIC.

Receiver: Intel X520 Server Adapter with Flow Director enabled (configured with suggested default parameters [11]: FdirMode=1, AtrSampleRate=20), 10Gbps, MTU 1500; Linux 2.6.34 (Multiple TX Queue Capable).

In our experiments, *iperf* was used to send n parallel TCP streams from sender to receiver for 100 seconds. *iperf* was not pinned to a specific core in the receiver. Linux was configured to run in *multicore peak performance* mode; the scheduler tries to use all core resources in parallel as much as possible, distributing the load equally among the cores. *Iperf* is a multi-threaded network application. With multiple parallel TCP data streams, a dedicated child thread is spawned and assigned to handle each stream. As a result, *iperf* threads may migrate across cores. The receiver was instrumented to record out-of-order packets, and we calculated relevant packet reordering ratios. The experiment results, with a 95% confidence interval, are shown in Table 1. The degree of packet reordering is significant. At $n=200$, packet reordering ratio reaches as high as 0.897%.

n	Reordering Ratio
40	0.498% \pm 0.067%
100	0.705% \pm 0.042%
200	0.897% \pm 0.038%
500	0.635% \pm 0.154%
1000	0.409% \pm 0.009%
2000	0.129% \pm 0.003%

Table 1 Experiment Results

The experiment results validated our analysis. When the scheduler tries to use all core resources in parallel as much as possible, distributing the load equally among the cores, it will lead to frequent process migration. As our analysis suggested, the Flow Director mechanism would cause packet reordering when process migration occurs. In addition, we ran `tcpdump` to record a single stream's packet trace at the receiver @ $n=200$. The packet trace analysis in Appendix B shows the occurrence of duplicate Acknowledgements (ACKs), SACKs, and data retransmissions due to packet reordering.

We then pin `iperf` to core 0 in the receiver and repeated the above experiments. No packet reordering was discovered. This is because when `iperf` is pinned to a specific core, its child threads are also pinned to that core. There will be no process migration in this case. In these conditions, Flow Director does not cause packet reordering.

The root cause of the packet reordering is that Flow Director lacks mechanisms to ensure in-order packet delivery when it steers packet across cores. In high-speed networks, packet reordering causes various negative impacts [8][9]. Many TCP implementations use the header prediction algorithm to reduce the costs of TCP processing. However, header prediction only works for in-sequence TCP segments. If segments are reordered, most TCP implementation do far more processing than they would for in-sequence delivery, degrading the TCP sender and receiver's performance. In addition, TCP SACK is now implemented and enabled by almost all general-purpose OSes. When packet reordering occurs, the receiver will sort the out-of-order queue to generate SACK blocks. For the sender, on receipt of SACK information, the retransmission queue would be walked and the relevant packets tagged as sacked or lost. In high-speed networks, the number of packets in the fly is large. The sender's retransmission queue is large. Also, when packet reordering occurs, out-of-order queue will be very large. Sorting out-of-sequence queue in the receiver or walking the retransmission queue in the sender can seriously degrade system performance [8][10]. Because the networking community is working towards 40GigE and 100GigE, the performance requirements on TCP/IP are becoming more challenging. Flow Director's packet reordering problem becomes more serious.

3. A Transport Friendly NIC (A-TFN)

3.1 A-TFN Design

We propose A-TFN mechanism to remedy the RSS and Flow Director limitations. A-TFN steers incoming network data to the same core on which its application thread resides and ensures in-order packet delivery. Our data steering mechanism is mainly targeted at TCP, but can be extended to UDP and SCTP. We base our A-TFN design on two observations. First, a TCP connection's traffic is bidirectional. For a unidirectional data flow, ACKs on the reverse path result in bidirectional traffic. Second, when an application makes socket-related system calls, that application's process context would be borrowed to carry out network processing in process context. This is true and common for all general purpose OSes although their network stacks are implemented differently. In

the data transmit path, network processing starts in the process context when an application makes socket-related system calls to send data. If TCP gives permission to send, network processing in process context can reach down to the bottom of the protocol stack. In the data receive path, when an application makes socket-related receive system calls to moves data from the socket into the user space, it needs to generate ACKs to advertise new receive window sizes. These ACKs are generated in process context.

A-TFN's basic idea is simple: it maintains the relationship "Traffic Flows \rightarrow Cores (Applications) in the NIC, with OSes correspondingly enhanced to support such capability. For transport layer traffic, A-TFN maintains a Flow-to-Core table in the NIC, with one entry per flow. Each entry tracks which receive queue (core) a flow should be assigned to. With each outgoing transport-layer packet (including ACK packet), the OS records a processor core ID and uses it to update the entry in the Flow-to-Core table. As soon as any network processing is performed in a process context, A-TFN learns of the core on which an application thread resides and can steer future incoming traffic to the right core. This is a key point that A-TFN is different from Flow Director.

The design of such a mechanism involves a trade-off between the amount of work done in the NIC and in the OS. There are two design options. Option 1 is to minimize changes in the OS and focuses instead on identifying the minimal set of mechanisms to add to the NIC. This design adds complexity and cost to the NIC. On the other end of the design space, it could be let the OS update the flow-to-core table directly without changing anything in the NIC hardware (option 2). Conceptually, this approach could be fairly straightforward to implement. However, it might add significant extra communication overheads between the OS and the NIC, especially when the Flow-to-Core table gets large. Due to space limitation, this paper is mainly focused on the first design option. In our future work, we will explore the second design option. Besides, option 1 design has other goals: (1) A-TFN must be simple and efficient. NIC controllers usually utilize a less powerful CPU with a simplified instruction set and insufficient memory to hold complex firmware. (2) A-TFN must preserve in-order packet delivery. (3) The communication overheads between the OS and A-TFN must be minimal.

Fig. 8 illustrates the A-TFN details. A-TFN extends the current RSS technologies. It supports multiple receive queues in the NIC, up to the number of cores in the system. With MSI and Flow-Pinning support, each receive queue has a dedicated interrupt and is tied to a specific core. Each core in the system is assigned a specific receive queue. A-TFN handles non-transport layer traffic in the same way as does RSS. That is, based on a hash of the incoming packet's headers, the NIC assigns it to the same queue as other packets from the same data flow, and distributes different flows across queues. For transport layer traffic, A-TFN maintains a Flow-to-Core table with a single entry per flow. Each entry tracks the receive queue (core) to which a flow should be assigned. The entries within the Flow-to-Core table are updated by outgoing packets. For unidirectional TCP data flows, outgoing ACKs update the Flow-to-Core table. For an outgoing transport-layer packet, the OS records a processing core ID in the transmit

descriptor and passes it to the NIC. Since each packet contains a complete identification of the flow it belongs to, the specific Flow \rightarrow Core relationship could be effectively extracted from the outgoing packet and its accompanying transmit descriptor. As soon as any network processing is performed in process context, A-TFN learns of on which core an application thread resides.

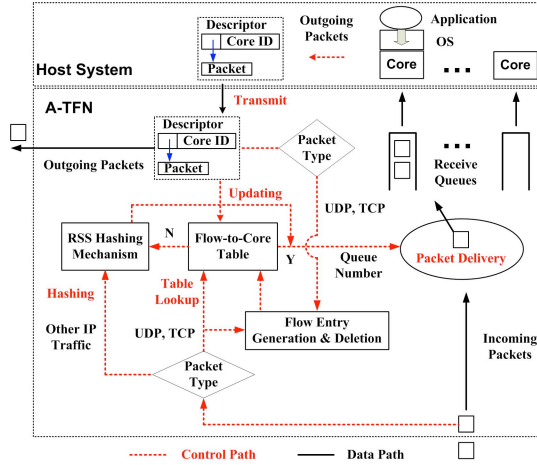


Fig. 8 A-TFN Mechanisms

3.2 Flow-to-Core Table and its Operations

The Flow-to-Core table consists of flow entries. Flow entries are managed in a hash table, with a linked list to resolve collisions. Each entry consists of:

- **Traffic Flow.** A-TFN makes use of the 5-tuple $\{src_addr, dst_addr, protocol, src_port, dst_port\}$ in the receive direction to specify a flow. Therefore, for an outgoing packet with the header $\{(src_addr: x), (dst_addr: y), (protocol: z), (src_port: p), (dst_port: q)\}$, its corresponding flow entry in the table is identified as $\{(src_addr: y), (dst_addr: x), (protocol: z), (src_port: q), (dst_port: p)\}$.
- **Core ID.** The core to which the flow should be steered.
- **Transition State.** A flag to indicate if the flow is in a transition state. The goal is to ensure in-order packet delivery.
- **Packets in Transition.** A simple packet list to accommodate temporary packets when the flow is in a transition state. The goal is to ensure in-order packet delivery.

In addition, to avoid non-deterministic packet processing time, a collision-resolving linked list is limited to a maximum size of $MaxListSize$. Flows are not evicted in case of collision. When a specific hash's collision-resolving list reaches $MaxListSize$, later flows with that hash will not be entered into the table.

a). Flow Entry Generation and Deletion

A-TFN monitors each incoming and outgoing packet to maintain the Flow-to-Core Table. An entry is generated in the Flow-to-Core table as soon as A-TFN detects a successful three-way handshake. To reduce NIC complexity, A-TFN need not run a full TCP state machine in the NIC. A flow

entry is deleted after a configurable period of time, T_{delete} , has elapsed without traffic. In this way, A-TFN need not handle all exceptions such as missing FIN packets and various timeouts. To prevent memory exhaustion or malicious attacks, A-TFN sets an upper bound on the number of entries in the Flow-to-Core Table. When the Flow-to-Core table starts to become full, TCP flows can be aged out more aggressively by using a smaller T_{delete} . For traffic flows that are not in the Flow-to-Core table, packets are delivered based on a hash of the incoming packets' headers.

b). Detection and Prevention of Packet Reordering

The entries of the Flow-to-Core table are updated by outgoing packets. For each outgoing transport-layer packet, the OS records a processing core ID in the transmit descriptor and passes it to the NIC. A naive way to update the corresponding flow entry is with the passed core ID, omitting any other measures. As soon as any network processing is performed in process context, A-TFN will learn of the process migration and can steer future incoming traffic to the right core. However, this simple flow entry updating mechanism cannot guarantee in-order packet delivery. In Section 2.5, we analyze why Flow Director cannot guarantee in-order packet delivery. The model and analysis can be also applied here. As we have analyzed, if ϵ is small and $n > m$, the condition of $T_{service}(S) > T_{service}(S+1)$ would easily hold and lead to packet reordering. Since the ring buffer size is D , the worst case is $n = D - 1$ and $m = 0$. It would have $T_{service}(S) = T - \epsilon + (D - 1) / R_{service}$ and $T_{service}(S + 1) = T + \epsilon$. TCP performance suffers in the event of severe packet reordering [8]. However, if the delivery of packet $S + 1$ to Core 1 can be delayed for at least $(D - 1) / R_{service}$, then $T_{service}(S + 1) \geq T + \epsilon + (D - 1) / R_{service}$. As a result, $T_{service}(S + 1) > T_{service}(S)$ and in-order packet delivery can be guaranteed. Therefore, A-TFN adopts the following flow entry updating mechanism: for each outgoing transport-layer packet, the OS records a processing core ID in the transmit descriptor and passes it to the NIC to update the corresponding flow entry. For a TCP flow entry, if the new core id is different from the old one, the flow enters the "transition" state. Correspondingly, its "Transition State" is set to "Yes" and a timer is started for this entry. The timer's expiration value is set to $T_{timer} = (D - 1) / R_{service}$. Incoming packets of a flow in the transition state are added to the tail of "Packets in Transition" instead of being immediately delivered. When the timer expires, the flow leaves the transition state. The "Transition State" is set back to "No" and all of the packets in "Packets in Transition," if they exist, are assigned to the new core. For a flow in the "non-transition" state, its packets are directly steered to the corresponding core. With current computing power, $(D - 1) / R_{service}$ is usually at the sub-millisecond level, at best. For A-TFN, T_{timer} is a design parameter and is configurable. In contrast, Flow Director does not have an effective mechanism to ensure in-order packet delivery.

3.3 Required OS Support

A-TFN design requires only two small OS changes in order to be properly supported. These can be easily implemented. (1) For an outgoing transport-layer packet, the OS needs to record a processing core ID in the transmit descriptor passed to the NIC. (2) The transmit descriptor needs to be updated with a new element to store this core ID. A single-byte element can support up to 256 cores, which is sufficient for most of today's systems. In addition, the size of a transmit descriptor is usually small, typically less than a cache line. Transmit descriptors are usually copied to the NIC by DMA using whole cache line memory transactions. Adding a byte to the transmit descriptor introduces almost no extra communication overhead between the OS and NIC.

4. Analysis and Experiments

The A-TFN mechanism is simple and requires the most minimal OS support. In addition, the communication overheads between the OS and A-TFN are reduced to a minimum. A-TFN can be effectively implemented with current hardware and software technologies.

4.1 Analytical Evaluation

a) Delay. To ensure in-order packet delivery, incoming packets of a flow in the transition state are added to the tail of “**Packets in Transition**”. These packets are delivered later, when the flow exits the transition state. Clearly, this can add delay to certain packets and the maximum delay a held packet can experience is T_{timer} . Previous analysis has shown that in-order packet delivery is guaranteed when T_{timer} is set to $(D-1)/R_p$. But incoming packets rarely fill a ring buffer in the real world. If T_{timer} were configured to be smaller, this would still ensure in-order packet delivery in most cases. We had recorded the duration for which the OS processes the ring buffer in [8]. The duration is generally shorter than 20 microseconds. In most cases the extra delay is so small that it can be ignored.

b) Flow Affinity and Network Data Affinity. The intent of A-TFN is to automatically steer incoming network data to the same core on which its application thread resides. As soon as any network processing is performed in a process context, A-TFN learns of the core on which an application thread resides and can steer future incoming traffic to the right core. The desired flow affinity and network data affinity are guaranteed.

c) Hardware design considerations. A-TFN's memory is mainly used to maintain the Flow-to-Core table, holding flow entries and accommodating packets for flows in the transition state. To hold a single flow entry, 20 bytes is quite sufficient. Therefore, a 10,000-entry Flow-to-Core table requires only 0.2 MB of memory. (These figures apply to IPv4; IPv6 support would add 24 bytes to the size of each entry, or less if the flow label could be relied upon.) In addition, to accommodate packets for flows in transition, if T_{timer} is set to 0.2 millisecond, even for a 10Gbps NIC, the memory required is $0.2\text{ ms} \times 10\text{Gbps} = 0.25\text{ MB}$, at maximum. In the worst case, an extra 0.5 MB of fast SRAM is enough to support the Flow-to-Core

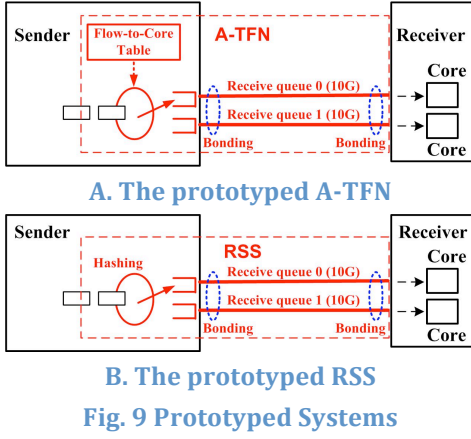
Table. A Cypress 4Mb (10ns) SRAM now costs around \$7. Appendix C lists the cost, memory size and power consumption of three popular 10G Ethernet NICs in the market. A-TFN's requirement of an extra 0.5 MB fast SRAM in the NIC won't add much extra cost and power consumption to current 10Gbps NICs. There is other hardware implementation cost. A-TFN might utilize content-addressable memories (CAMs) to implement the lookup function in the flow-to-core table. A linked list in HW is expensive to build given all the extra handling. There will be a tradeoff in hardware complexity (cost) and A-TFN effectiveness.

d) 40GigE and 100GigE. The networking community is working towards 40GigE and 100GigE. A-TFN must be applicable to these emerging technologies. In fact, A-TFN can be simply extended to 40GigE or 100GigE except a few small changes. First, *MaxListSize*, the maximum size of the collision-resolving linked lists of the Flow-to-Core table, should be reduced if higher memory speed is not available. For a 10GbE NIC, the time budget to process a 1500byte packet is around 1200 ns. For a 40GigE and 100GigE NIC, such time budget is reduced to 300 ns and 120 ns, respectively. In reality, A-TFN's actual allowable time budget to process a packet is even smaller due to the existence of smaller sized packets (<1500bytes). Assume A-TFN's other operations such as hash computing and packet delivery totally take T_{other} and each item in a linked list takes an extra T_{item} to access. Therefore, the *MaxListSize* for a 40GigE NIC and 100GigE NIC is approximately $(300 - T_{other})/T_{item}$ and $(120 - T_{other})/T_{item}$, respectively. Second, more memory is required to hold packets for flows in transition to ensure in-order packet delivery. If T_{timer} is set to 0.2 millisecond, for a 40GigE NIC, the memory required is $0.2\text{ ms} \times 40\text{Gbps} = 1\text{ MB}$; for a 100GigE NIC, the memory required is 2.5 MB.

4.2 Experimental Evaluation

4.2.1 Prototyped Systems

We prototyped an A-TFN with two receive queues as shown in Fig. 9A. A sender connects to a receiver via two physical back-to-back 10Gbps links. The sender and receiver are the same computer systems as specified in Section 2.4. The 10Gbps links are driven by Myricom 10Gbps Ethernet NICs. In both the sender and the receiver, the two Myricom 10Gbps NICs are aggregated into a single logical bonded interface with the Linux bonding driver. In the sender, the bonding driver is modified with A-TFN mechanisms and each 10Gbps link is deemed an A-TFN receive queue. In the receiver, each slave NIC (receive queue) is pinned to a specific core. In addition, the receiver's OS is modified to support the A-TFN mechanisms. For an outgoing transport-layer packet, the OS records a processing core ID in the “transmit descriptor” and passes it to “A-TFN.” Here, we make use of four reserved bits in the TCP header as the “transmit descriptor” to communicate the core ID. When the sender receives a “transmit descriptor,” it extracts the passed Core ID and updates the corresponding flow entry in the Flow-to-Core table. Unless otherwise



specified, T_{timer} is set to 0.1 ms. The Flow-to-Core table is upped limited to 10,000 entries. In our emulated system, we measure the Flow-to-Core Table's search time. The search time to access the first item in a collision-resolving linked list takes around 260 ns, which includes the hashing and locking overheads. For each next item in the list, it takes approximately an extra 150 ns. Therefore, the longest search in our system takes $260 + 150 * (MaxListSize - 1)$ ns. For a 10Gbps NIC, the time budget to process a 1500byte packet is around 1200 ns. To evaluate $MaxListSize$'s effect on A-TFN's performance, we set $MaxListSize$ to 1 and 6, respectively. Correspondingly, A-TFN is termed as A-TFN-1 and A-TFN-6.

Similarly, we implemented a two-receive queue RSS NIC, as shown in Fig. 9B. In both the sender and the receiver, the two Myricom 10Gbps NICs are aggregated into a single logical bonded interface with the bonding driver. In the sender, the bonding driver is modified with RSS mechanisms, and each 10Gbps link is treated as an RSS receive queue. Unless otherwise specified, the hashing is based on the combination of $\{src_addr, dst_addr, src_port, dst_port\}$ for each incoming packet. In the receiver, each slave NIC (receive queue) is pinned to a specific core.

4.2.2 Experiment Configurations

We ran data transmission experiments with iperf using the prototyped systems shown in Fig. 9. In our experiments, iperf sends with n parallel TCP streams for 100 seconds, to ports 5001 and 6001, respectively. Therefore, totally $2n$ parallel TCP streams are transmitting in each experiment. The number n was varied across experiments. In all the experiments, the sender runs the same scripts. The scrip runs in the sender is:

```
iperf -c receiver -P n -t 100 -p 5001 &
iperf -c receiver -P n -t 100 -p 6001 &
```

The experiment configurations in the receiver are varied across experiments.

Experiment 1 was designed to verify that A-TFN can remedy the RSS limitation. In section 2.2 we discussed four cases in OSES in which a network application thread resides on cores other than those to which its corresponding network interrupts are directed. Due to page limitation, we only discussed the case that a single-threaded application must

handle multiple concurrent TCP connections in the experiment. In the real world, there are many cases that a single-threaded application must handle multiple concurrent TCP connections. For example, Nginx [22] and Lighttpd [23] are such cases. Nginx and Lighttpd are probably the two best-known asynchronous HTTP servers. They are event-driven and handle multiple concurrent TCP connections in a single thread (or at least, very few threads). In Experiment 1, TCP streams of a specific port (5001 or 6001) were pinned to a particular core in the receiver (Table 2). In this way, we simulated iperf of port 5001 and 6001 as two "single-threaded" applications that run on core 0 and 2, respectively. Each single-threaded application handles n concurrent TCP connections. There are a few reasons why we simulated a single-threaded application using iperf (a multi-threaded application). First, the purpose of the experiments is to verify that A-TFN can effectively improve network performance and Iperf is a simple and commonly used networking testing tool. Second, if a real single-threaded application like Nginx were used in the experiments, the complicated software itself might interfere with the network performance testing. For example, Nginx is event-driven and many of its activities are unrelated to network operations.

Receive Queues	Iperf Config.
ReceiveQ 0 @ Core 0	"iperf -s -p 5001" @ Core {0}
ReceiveQ 1 @ Core 2	"iperf -s -p 6001" @ Core {2}

Table 2 Experiment 1 Receiver Configurations

Different from Flow Director, A-TFN uses a special flow entry updating mechanism to guarantee in-order packet delivery. Experiment 2 was designed to evaluate whether this mechanism actually works. In Experiment 2, iperfs (ports 5001 and 6001) were allowed to run on both cores where the two receive queues were pinned (Core 0 and 2) (Table 3). Linux was configured to run in *multicore peak performance* mode. As a result, iperf threads may migrate across cores.

Receive Queues	Iperf Config.
ReceiveQ 0 @ Core 0	"iperf -s -p 5001" @ Core {0, 2}
ReceiveQ 1 @ Core 2	"iperf -s -p 6001" @ Core {0, 2}

Table 3 Experiment 2 Receiver Configurations

4.2.3 Experiment Results

a) Experiment 1 Results

Given the same experimental conditions, we compared the results with A-TFN to those with RSS. The metrics of interest were: (1) Throughput; (2) WAITTIME-TOTAL; and (3) BUS_HITM_DRV. (The number of events between samples was 10000.) Consistent results were obtained across repeated runs. All results presented are shown with a 95% confidence interval.

When a single-threaded network application handles multiple concurrent TCP connections, the hashing function of the RSS-enabled NIC will evenly and statistically distribute the connections across the cores. Since the application can

only run on a single core at any given moment, some connections get steered to cores other than the one on which the application runs. As a result, TCP processing will alternate between different cores. This fact may even lead to contention for shared resources between interrupt and application threads when they concurrently process packets of the same flows. Under such circumstances, overall system efficiency could be severely degraded. The experimental results in Fig. 10 confirm these points. Experiment 1 shows that: (1) A-TFN can effectively improve the network throughput. A-TFN-6 markedly increased the TCP throughput by more than 20% with $2n=1000$. (2) A-TFN can significantly reduce lock contention in parallel network stacks. The total time spent waiting to acquire various kernel locks was decreased by more than 98% for A-TFN-6 with $2n=40$. (3) A-TFN can substantially reduce system synchronization overhead. Experimental data confirms the effectiveness of A-TFN in improving network throughput and enhancing system efficiency. This is because the design of A-TFN steers incoming network traffic to the same core(s) on which its application thread resides. Therefore, TCP processing does not alternate between different cores and contention involving shared resources between interrupt and application threads will not occur. In addition, costly MESI operations can be greatly reduced. For Experiment 1, the improvements in synchronization and cache statistics are substantial, yet they do not seem result in equivalent gains in throughput. This is because TCP data transmission involves complex interaction of the sender and receiver. Certainly, the improvements in synchronization and cache statistics in the receiver only cannot result in equivalent gains in TCP throughput. In the

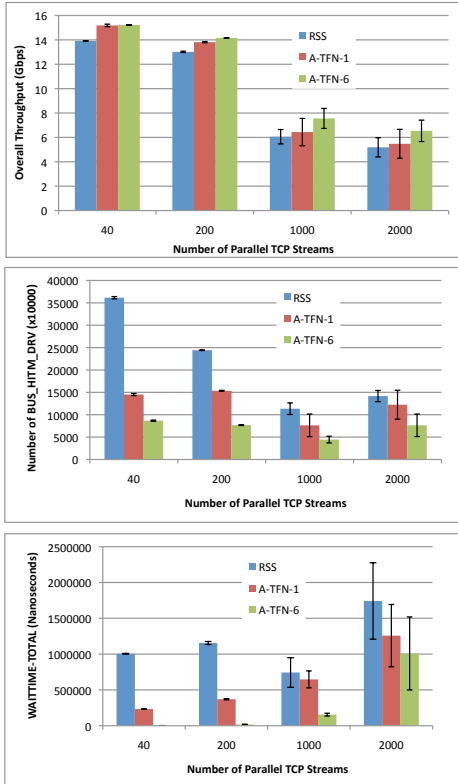


Fig. 10 Experiment 1 Results

experiments we also noticed that the receiver's system bus approaches saturation, which also partially explains the phenomena. We believe that if faster system bus were applied, the improvements in synchronization and cache statistics in the receiver would result in more gains in throughput.

For the Flow-to-Core table, when a specific hash's collision-resolving lined list reaches *MaxListSize*, subsequent flows for that hash will not be entered into the table. Their packets are delivered in the same way as does RSS. It can be seen from Fig. 10 that with $2n=40$, A-TFN-1's results (especially for throughputs) are very close to those of A-TFN-6's. With $2n=2000$, A-TFN-1 behaves closer as does RSS. We record the percentage of flows that are entered into the Flow-to-Core table when n is varied in Table 4. It shows that as n increases, the percentage of flows that are entered into the table decreases, with the effects on A-TFN-1 being much more than on A-TFN-6. With $2n=2000$, A-TFN-1 has only a 12.7% of flows entered into the Flow-to-Core table. The reason the ratio is so low is because all the flows share a single pair of IP addresses, they are not hashed efficiently across the table. As a result, more traffic would be delivered in the same way as RSS does. From the hardware implementation's perspective, A-TFN-1's Flow-to-Core table is much easier to implement. But its performance is not satisfactory as the number of TCP streams increase. Thus, there will be a tradeoff in hardware complexity (cost) and A-TFN effectiveness. It is anticipated that with n further increased, A-TFN-6 would have more traffic delivered in the way as RSS does; its effectiveness would start to decrease as well. Normally, a high-end web server would handle a few thousand concurrent TCP streams. For our two-core A-TFN emulated system, 2000 streams is quite a high number. Since the trend is already very clear, we don't further increase n .

$2n$	A-TFN-6	A-TFN-1
40	100% \pm 0	88% \pm 1.6%
200	100% \pm 0	71% \pm 2.9%
1000	95.7% \pm 1.1%	24.5% \pm 0.1%
2000	71.7% \pm 0.2%	12.7% \pm 0%

Table 4 Flows @ Flow-to-Core Table Percentage

With RSS technologies, the worst cases occur when soft partition technologies, like CPuset, are applied in the networking environments. This can easily lead to the undesirable situation in which network applications are soft-partitioned on cores other than those to which their network interrupts are directed. Also, an OS scheduler prioritizes load balancing (or power saving) over process-to-interrupt affinity. In these environments, network applications may also be scheduled on cores other than those where their corresponding network interrupts are directed. We ran experiments in these environments. All the experiments verify that A-TFN can steer incoming network data to the same core on which its application thread resides, resulting in improved performance.

b) Experiment 2 Results

In Experiment 2, iperfs were allowed to run on both cores and Linux was configured to run in *multicore peak*

performance mode. Therefore, iperf threads may migrate across cores. The receiver was instrumented to record out-of-order packets and we calculated relevant packet reordering ratios. For A-TFN-6, we set T_{timer} to 0 and 100 μ s, respectively. The experimental results are shown in Table 5.

When T_{timer} is 0, incoming packets of a flow in the transition state are immediately delivered, instead of being added to the tail of “*Packets in Transition*.” As discussed before, this could lead to packet reordering. The results in Table 5 reflect this fact. However, it can be seen that with $2n=40$ and $2n=200$, the packet reordering ratio is pretty low. This is because our experiments were actually run in a two-core system. When n is low, fewer iperf threads are spawned and process migration would occur less frequently. Therefore, less packet reordering would occur. We believe that if the experiments were run in a system with more cores, process migration would occur more frequently and would lead to more incidences of packet reordering even if n were low. On the other hand, it can be seen that when n is further increased (with $2n=1000$ and $2n=2000$), the packet ratio steadily increases. This is because when n is increased, more iperf threads are spawned and process migration will occur more frequently in the simulated two-core system. As a result, more packet reordering will result.

When T_{timer} is 100 μ s, no out-of-order packets are recorded. This shows that A-TFN’s packet reordering prevention mechanism really takes effect and can effectively guarantee in-order packet delivery.

$2n$	$T_{timer} = 0 (\mu s)$	$T_{timer} = 100 (\mu s)$
40	$5.110E-07 \pm 6.809E-07$	0
200	$6.278E-06 \pm 8.553E-06$	0
1000	$3.639E-05 \pm 2.754E-05$	0
2000	$2.174E-04 \pm 8.515E-05$	0

Table 5 Packet Reordering Ratios

5. Related Works

Over the years, research on affinity in network processing has been extensive. Salehi et al. [2] studied the effectiveness of affinity-based scheduling in multiprocessor network protocol processing using both packet-level and connection-level parallelization approaches. But since these approaches worked in the user space, they did not consider either system or implementation costs. A. Foong et al. [3] experimented with affinizing processes/threads, as well as interrupts from NICs, to specific processors in an SMP system. Experimental results suggested that processor affinity in network processing contexts can significantly improve overall performance. J. Hye-Churn et al. [4] studied the problem of multi-core aware processor affinity for TCP/IP over multiple network interfaces, using a software-only approach. Their research topics are similar to us.

Other researchers have adopted a hard partition approach [24][25]. In multiprocessor environments, a subset of the processor is dedicated to network processing; the remaining processors perform only

application-relevant computations. The limitation of this approach is that the OS architecture requires significant changes.

The NIC technologies, such as Intel’s vmdq [26] or the PCI-SIG’s SR-IOV [27], also provide data steering capabilities for the NICs. But they are I/O virtualization technologies targeting at virtual machines in the virtualized environment, not targeting at general purpose OSes in the non-virtualized environment. Intel Ethernet Flow Director technology [7] can automatically steer incoming network data to the same core on which its application thread resides. However, Flow Director can cause significant packet reordering in multiprocessing environments.

The Receive Packet Steering (RPS) [28] and Receive Flow Steering (RFS) [29] technologies are recently introduced. Both RPS and RFS are OS software technologies, instead of NIC technologies. They make use of an extra core in a multicore system to spread and steer incoming packets to other cores. RPS and RFS complement the RSS and A-TFN mechanisms. They are applied when NIC does not support RSS or A-TFN.

6. Conclusion and Discussion

We propose an A-TFN mechanism to remedy the limitations in RSS and Flow Director. In the paper, we discuss two A-TFN design options. Due to space limitation, this paper is mainly focused on the first design option. The new NIC is emulated in software. The experimental results show our solution is effective and practical to remedy the limitations we have identified in RSS and Flow Director. In future work, we will explore the second design option.

In our experiments, the sender and receiver are connected back-to-back. As a result, the Round Trip Time (RTT) is less than 0.1ms. With such a small RTT, the packet reordering’s negative impacts cannot take full effect. That is the reason why we did not present experimental evidence of the impacts of various degree of reordering on the overall performance in the paper. Luckily, there are various previous researches that studied the packet reordering’s negative impacts, which are convincing. And we have cited these researches in the paper. Readers might ask why not we run the experiments with larger RTTs? The answer is simple: we cannot run such experiments due to Limits of Current Experiment Conditions. The maximum throughput in our experiments is close to 15Gbps. In the real world, it is difficult for us to find a network path with RTT at least greater than 5ms and with bandwidth greater than 15Gbps to run our experiments. Fermilab does have such networking facilities to other sites. But these networks run production traffic. We are not allowed to run such experiments in our production networks. Very few people in the world, if not none, can find suitable networks to run our experiments. Furthermore, we cannot emulate such a network path in the lab environments either. There are tools like Netem [30] that provides network emulation functionality by emulating the properties of wide area networks. However, almost all these tools do no work well in high-speed networks (>5Gbps) due to system clock resolution issues or system bus speed issues.

References:

-
- [1] P. Willmann et al., "An Evaluation of Network Stack Parallelization Strategies in Modern Operating Systems," In Proc. USENIX Annual Technical Conference, pp. 91–96, 2006
- [2] J. D. Salehi, J. F. Kurose, and D. Towsley, "The effectiveness of affinity-based scheduling in multiprocessor network protocol processing (extended version)," *IEEE/ACM Trans. Networking*, vol. 4, pp. 516–530, Aug. 1996.
- [3] A. Foong et al., "An in-depth analysis of the impact of processor affinity on network performance," In Proc. IEEE International Conference on Networks, 2004.
- [4] J. Hye-Churn et al., "MiAMI: Multi-Core Aware Processor Affinity for TCP/IP over Multiple Network Interfaces," In Proc. IEEE Symposium on High Performance Interconnects, 2009.
- [5] R. Huggahalli et al., "Direct Cache Access for High Bandwidth Network I/O," In Proc. 32nd Annual International Symposium on Computer Architecture, 2005.
- [6] Microsoft Corporation (Nov. 2008), Receive-side scaling enhancements in windows server 2008. [Online]. Available: http://download.microsoft.com/download/a/d/f/adf1347d-08dc-41a4-9084-623b1194d4b2/RSS_Server2008.docx
- [7] Intel Corporation (Nov. 2010), Intel 82599 10gbe controller datasheet. [Online]. Available: http://download.intel.com/design/network/datashts/82599_datasheet.pdf
- [8] W. Wu et al., "Sorting reordered packets with interrupt coalescing," computer network, Volume 53, Issue 15, 2009, pages: 2646-2662.
- [9] M. Laor, L. Gendel, "The effect of packet reordering in a backbone link on application throughput," *IEEE Network*, vol. 16, no. 5, pp. 28-36, Sep/Oct 200.
- [10] P. McManus, Performance tradeoffs of TCP Selective Acknowledgement. [Online]. Available: <http://www.ibm.com/developerworks/linux/library/l-tcp-sack/>
- [11] — (Nov. 2010), Ixgbe device driver readme. [Online]. Available: <http://downloadmirror.intel.com/14687/eng/README.txt>
- [12] Linux Kernel Website. [Online]. Available: <http://kernel.org/>
- [13] S. Siddha et al., Chip Multi Processing aware Linux Kernel Scheduler, In Proc. the Linux Symposium, pp. 329 – 340, 2006
- [14] V. Pallipadi et al., "Processor Power Management features and Process Scheduler: Do we need to tie them together?" In Proc. LinuxConf Europe, 2007.
- [15] M.E. Russinovich et al., *Microsoft Windows Internals*, fourth ed., Microsoft Press, 2004. ISBN 0735619174.
- [16] W. Wu et al., The performance analysis of Linux networking – packet receiving, *Computer Communications* 30 (2007) 1044–1057.
- [17] W. Wu et al., Potential performance bottleneck in Linux TCP, *International Journal of Communication Systems* 20 (11) (2007) 1263–1283.
- [18] L. Ivanov et al., Modeling and verification of cache coherence protocols, In Proc. the IEEE International Symposium on Circuits and Systems, pp. 129 – 132, 2001
- [19] Iperf Website. [Online]. Available: <http://iperf.sourceforge.net/>
- [20] Oprofile Website. [Online]. Available: <http://oprofile.sourceforge.net/>
- [21] Intel 64 and IA-32 Architectures Software Developer's Manual, Volume 3B: System Programming Guide, Intel Corporation, 2008
- [22] NGINX Website. [Online]. Available: <http://nginx.org>
- [23] LIGHTTPD Website. [Online]. Available: <http://lighttpd.net>
-
- [24] T. Brecht et al. Evaluating Network Processing Efficiency with Processor Partitioning and Asynchronous I/O. In Proc. EuroSys, 2006.
- [25] G. Regnier et al., ETA: Experience with an Intel Xeon processor as a packet processing engine, *IEEE Micro*, 2004.
- [26] Intel Corporation, "Intel VMDq Technology," 2008.
- [27] <http://www.pcisig.com/specifications/iov>
- [28] <http://lwn.net/Articles/328339/>
- [29] <http://lwn.net/Articles/382428/>
- [30] Netem Website. [Online]. Available: <http://www.linuxfoundation.org/collaborate/workgroups/networking/netem>